

online softmax from math to code

wenqingqian May 2024

code see [Github](#)

Contents

1	Math: from 3 pass safe softmax to 2 pass online softmax	1
2	Code implementation	3
3	Profiling	5

1 Math: from 3 pass safe softmax to 2 pass online softmax

3 pass safe softmax

for i in $[1, N]$:

$$m_i = \max(x_i, m_{i-1}) \quad (1)$$

for i in $[1, N]$:

$$d_i = d_{i-1} + e^{x_i - m_N} \quad (2)$$

for i in $[1, N]$:

$$y_i = \frac{e^{x_i - m_N}}{d_N} \quad (3)$$

The equation 2 equals to:

$$d_N = \sum_{i=1}^N e^{x_i - m_N} \quad (4)$$

which comes from:

$$d_j = \sum_{i=1}^j e^{x_i - m_N} \quad (5)$$

Actually we can create another sequence to remove the dependency on N :

$$d'_j = \sum_{i=1}^j e^{x_i - m_j} \quad (6)$$

when $j = N$, $d_j = d'_j$.

Then we can find the recurrence relation between d'_j and d'_{j-1} :

$$\begin{aligned} d'_j &= \sum_{i=1}^j e^{x_i - m_j} \\ &= \sum_{i=1}^{j-1} e^{x_i - m_j} + e^{x_j - m_j} \\ &= \sum_{i=1}^{j-1} e^{x_i - m_{j-1} + m_{j-1} - m_j} + e^{x_j - m_j} \\ &= \sum_{i=1}^{j-1} (e^{x_i - m_{j-1}} * e^{m_{j-1} - m_j}) + e^{x_j - m_j} \\ &= \left(\sum_{i=1}^{j-1} e^{x_i - m_{j-1}} \right) * e^{m_{j-1} - m_j} + e^{x_j - m_j} \\ &= d'_{j-1} * e^{m_{j-1} - m_j} + e^{x_j - m_j} \end{aligned} \quad (7)$$

2 pass online softmax

for i in $[1, N]$:

$$\begin{aligned} m_i &= \max(x_i, m_{i-1}) \\ d'_i &= d'_{i-1} * e^{m_{i-1} - m_i} + e^{x_i - m_i} \end{aligned} \quad (8)$$

for i in $[1, N]$:

$$y_i = \frac{e^{x_i - m_N}}{d'_N} \quad (9)$$

2 Code implementation

This algorithm is easy to implement in CPU, but in GPU, to make the algorithm parallel, we need to generalize the original algorithm and use some CUDA tools.

In equation 8, each time we process a new number x_i , and calculate the new max number of sequence and summation. To make the algorithm parallel, we need to get the new result by two sequence directly, assuming there are two sequences:

- S with
length: l_S
max: $S.\text{max}$
summation: $S.\text{sum} = \sum_{i=1}^{l_S} e^{s_i - S.\text{max}}$
- B with
length: l_B
max: $B.\text{max}$
summation: $B.\text{sum} = \sum_{i=1}^{l_B} e^{b_i - B.\text{max}}$

The sequence $A = S + B$, assuming $S.\text{max}$ is smaller than $B.\text{max}$.

- A
length: $l_S + l_B$
max: $A.\text{max} = B.\text{max}$
summation: $A.\text{sum} = \sum_{i=1}^{l_S+l_B} e^{a_i - A.\text{max}}$

Then the summation of A equals to:

$$\begin{aligned}
A.\text{sum} &= \sum_{i=1}^{l_S+l_B} e^{a_i-A.\text{max}} \\
&= \sum_{i=1}^{l_S+l_B} e^{a_i-B.\text{max}} \\
&= \sum_{i=1}^{l_S} e^{s_i-B.\text{max}} + \sum_{i=1}^{l_B} e^{b_i-B.\text{max}} \\
&= \sum_{i=1}^{l_S} e^{s_i-S.\text{max}+S.\text{max}-B.\text{max}} + \sum_{i=1}^{l_B} e^{b_i-B.\text{max}} \\
&= \left(\sum_{i=1}^{l_S} e^{s_i-S.\text{max}} \right) * e^{S.\text{max}-B.\text{max}} + \sum_{i=1}^{l_B} e^{b_i-B.\text{max}} \\
&= S.\text{sum} \cdot e^{S.\text{max}-B.\text{max}} + B.\text{sum}
\end{aligned} \tag{10}$$

Generally, we do not assume that $B.\text{max}$ is greater than $S.\text{max}$:

$$\begin{aligned}
A.\text{max} &= \max(B.\text{max}, S.\text{max}) \\
A.\text{sum} &= S.\text{sum} \cdot e^{S.\text{max}-A.\text{max}} + B.\text{sum} \cdot e^{B.\text{max}-A.\text{max}}
\end{aligned} \tag{11}$$

Then we can process the sequence by multiple threads through warp-reduce.

online softmax CUDA pseudo

Assuming *input* sequence length $N = 32$

-
- 1: In each thread $i \in [1, 32]$
 - 2: local max = *input*[i]
 - 3: local sum = $e^0 = 1$
 - 4: **for** (mask = $32 \gg 1$; mask ≥ 1 ; mask $\gg = 1$) **do**
 - 5: sum_tmp = __shfl_xor_sync(0xffffffff, local sum, mask)
 - 6: max_tmp = __shfl_xor_sync(0xffffffff, local max, mask)
 - 7: newmax = max(local max, max_tmp)
 - 8: local sum = local sum $\cdot e^{\text{local max} - \text{newmax}}$ + sum_tmp $\cdot e^{\text{max_tmp} - \text{newmax}}$
 - 9: local max = newmax

```

10: end for
11: In each thread  $i \in [1, 32]$ 
12:  $output[i] = (e^{input[i] - local\ max}) / (local\ sum)$ 

```

3 Profiling

As shown in the fig1, when the sequence size is larger than 512*16, the performance of online softmax is worse than 3 pass softmax. Since it is difficult to synchronize sum and max information between blocks, our algorithm only uses one block and 128 threads. When the data increases, multiple data will be allocated to one thread, and the processing of these data is serial, causing the time consumption to increase exponentially.

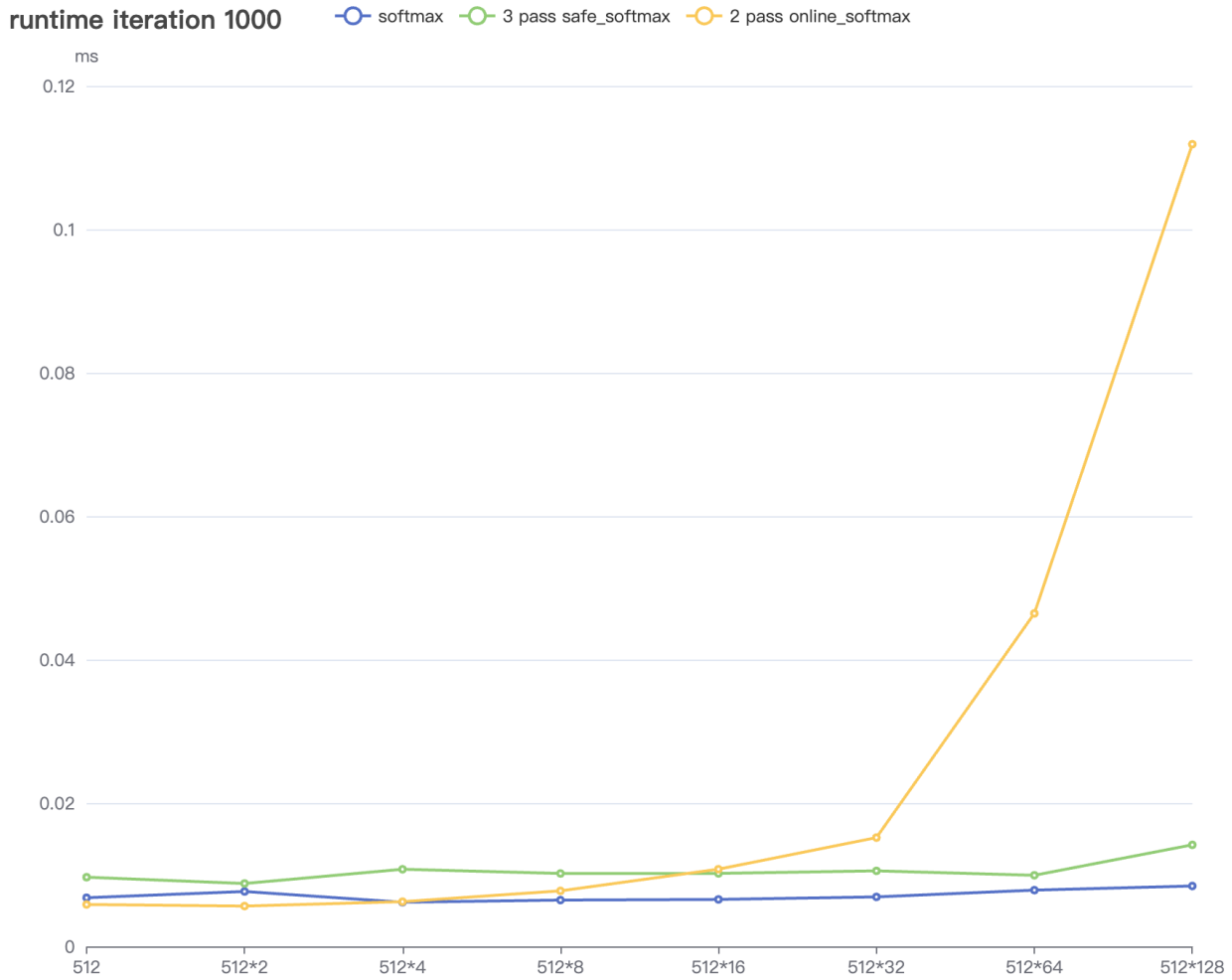


Figure 1: Profiling